

SECUR

AD-A233 649

ENTATION PAGE

DTIC FILE COPY

Form Approved
OMB No. 0704-0188

1a. RI	1b. RESTRICTIVE MARKINGS		
2a. SECURITY CLASSIFICATION AUTHORITY	3. DISTRIBUTION/AVAILABILITY OF REPORT		
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE	Approved for public release; distribution unlimited.		
4. PERFORMING ORGANIZATION REPORT NUMBER(S)	5. MONITORING ORGANIZATION REPORT NUMBER(S)		
6a. NAME OF PERFORMING ORGANIZATION The Regents of the University of California	6b. OFFICE SYMBOL (If applicable)	7a. NAME OF MONITORING ORGANIZATION Cognitive Science Program Office of Naval Research (Code 1142PT)	
6c. ADDRESS (City, State, and ZIP Code) University of California, Los Angeles Office of Contracts and Grants Administration Los Angeles, California 90024	7b. ADDRESS (City, State, and ZIP Code) 800 North Quincy Street Arlington, VA 22217-5000		
8a. NAME OF FUNDING/SPONSORING ORGANIZATION Defense Advanced Research Projects Agency	8b. OFFICE SYMBOL (If applicable)	9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER N00014-86-K-0395	
8c. ADDRESS (City, State, and ZIP Code) 1400 Wilson Boulevard Arlington, VA 22209-2308	10. SOURCE OF FUNDING NUMBERS		
	PROGRAM ELEMENT NO. 61153N	PROJECT NO. RR04206	TASK NO. RR04206-OC
	WORK UNIT ACCESSION NO. 442c022		
11. TITLE (Include Security Classification) A Sourcebook Approach to Evaluating Artificial Intelligence Systems			
12. PERSONAL AUTHOR(S) Dyer, Michael and Read, Walter			
13a. TYPE OF REPORT Interim	13b. TIME COVERED FROM 7/1/86 TO 4/30/88	14. DATE OF REPORT (Year, Month, Day) April 1988	15. PAGE COUNT 6
16. SUPPLEMENTARY NOTATION Paper presented at the Annual Meeting of the American Educational Research Association for the Symposium "How Smart are Smart Computers? Alternative Approaches to the Evaluation of Artificial Intelligence Technology," New Orleans, April 17.			
17. COSATI CODES 18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) 1988 Artificial intelligence, natural language processing, linguistics			
19. ABSTRACT (Continue on reverse if necessary and identify by block number) This paper provides the rationale and initial plan for the development of a sourcebook of natural language processing problems one of the research tasks under the natural language component of the Artificial Intelligence Measurement System (AIMS). The paper provides a prototype example of a sourcebook entry called an "exemplar."			
20. DISTRIBUTION/AVAILABILITY OF ABSTRACT <input checked="" type="checkbox"/> UNCLASSIFIED/UNLIMITED <input type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS		21. ABSTRACT SECURITY CLASSIFICATION Unclassified	
22a. NAME OF RESPONSIBLE INDIVIDUAL Dr. Susan Chipman		22b. TELEPHONE (Include Area Code) (703) 696-4318	22c. OFFICE SYMBOL ONR 1142CS

Project Report #7

A SOURCEBOOK APPROACH TO
EVALUATING ARTIFICIAL INTELLIGENCE SYSTEMS

Michael Dyer

Walter Read

Artificial Intelligence Laboratory
Computer Science Department
UCLA

This paper was presented at the April 1988 annual meeting of the
American Educational Research Association, New Orleans.

Artificial Intelligence Measurement System
Contract Number N00014-86-K-0395

Principal Investigator: Eva L. Baker

Center for Technology Assessment
UCLA Center for the Study of Evaluation

Acquisition File	
NTIS	Classification
UDC	Code
Document No.	Date Received
Justification	
By	
Distribution	
Administrative Control	
Dist	Approved by C.S.E. Special
A-1	

This research report was supported by contract number N00014-86-K-0395 from the Defense Advanced Research Projects Agency (DARPA), administered by the Office of Naval Research (ONR), to the UCLA Center for the Study of Evaluation. However, the opinions expressed do not necessarily reflect the positions of DARPA or ONR, and no official endorsement by either organization should be inferred. Reproduction in whole or part is permitted for any purpose of the United States Government.

Approved for public release; distribution unlimited.

Evaluating Natural Language Systems

Recent years have seen a proliferation of computer systems for natural language processing (NLP). These include front ends to databases, expert systems and tutoring systems. Such systems generally come with a list of inputs (typically single sentences) that the system is claimed to 'handle'. The problem in judging these systems is that it is very difficult to tell from the examples just what claims are being made. If one of the examples includes an ellipsis, does that mean the system handles ellipsis in general? Or only certain kinds? What is ellipsis 'in general'? Are there different kinds of ellipsis that require different kinds of understanding?

Evaluating these claims requires that we know what inputs the system *should* handle and what it would mean to understand the input. Testing understanding is easier for applied systems since there is generally a specific task involved, e.g., accessing a database. But deciding what inputs should be handled is more difficult because there is no general agreement on what kinds of linguistic phenomena there are. Without a common classification of the problems in natural language understanding authors have no way to specify clearly what their systems do, potential users have no way to compare different systems and researchers have no way to judge the advantages or disadvantages of different approaches to developing NLP systems.

This paper reports progress in development of evaluation methodologies for natural language systems. This work is part of the Artificial Intelligence Measurement System (AIMS) project of the Center for the Study of Evaluation at UCLA.

Previous Work

These problems have been discussed for some time in computer science NLP work but there has been very little work in developing actual evaluative criteria. Woods (1977) discussed the taxonomic approach and pointed out some of its strengths and weaknesses. Guida and Mauri (1984, 1986) discuss a formal model which involves measuring the correctness of the understanding and averaging it over a weighted set of inputs. But this method assumes that we can describe a weighting for (categories of) inputs.

The Sourcebook

In developing evaluative criteria for NLP systems we had several goals in mind. First, the criteria used should be applicable over the broadest possible range of systems and still provide comparability of the systems. Second, the system shouldn't just be rated on a pass/fail count. It should outline areas of competence so that implementers can see where further work is needed in their system. They should be able to say "this approach handles types 1, 2 and 3 of ellipsis but not types 4 and 5 yet" rather than "this approach handles ellipsis". Third, the criteria used should be comprehensible to the general user and to researchers outside computational linguistics. We need to present the issues in such a way that the user can make judgments about the importance of different components of the evaluation. This means presenting the issues in terms of the general principles involved and giving concrete examples. This approach also allows us to bring in information from areas like education, psychology, sociology, law and literary analysis and enables researchers in those areas to contribute to the evaluation.

To this end, we are building a database of *exemplars* of representative problems in natural language understanding, mostly from the computational linguistics literature. Each exemplar includes a piece of text (sentence, dialogue fragment, etc.) a description of the conceptual issue represented, a detailed discussion of the problems in understanding the text and a reference to a more extensive discussion in the literature. The Sourcebook consists of a large set of these exemplars and a conceptual taxonomy of the types of issues represented in the database. The exemplars are indexed by source in the literature and by conceptual class of the issue so that the user can readily access the relevant examples. The Sourcebook provides a structured representation of the coverage that can be expected of a natural language system.

Rather than start with a particular theory of language, we began with a search of the computational linguistics literature. While no-one would claim that computational linguistics has discovered, let alone solved, every problem in language use, twenty-five years of research has covered a broad range of problems. Looking at language use computationally focuses attention on phenomena that are often neglected in more theoretical analyses. Building systems intended to read real text or interact with real users raises complex problems of interaction of linguistic phenomena. The exemplars are mostly taken from the literature although we have added examples to fill in gaps

where we felt the published examples were incomplete. Because many of the published cases involved particular systems, the examples are often discussed in the literature in relation to that system. In the exemplars, we analyze the example in terms of the general issue represented. Then the exemplars are grouped into categories of related problems. This will generate the hierarchical classification of the issues.

Continuing and Future Work

We have several hundred exemplars and we estimate that we have covered 10 per cent of the relevant literature (journals, proceedings volumes, dissertations, major textbooks) in computational linguistics, artificial intelligence and cognitive science.

We are continuing to add exemplars to the Sourcebook and are elaborating the classification scheme. We will be making the Sourcebook available to other researchers for comment and analysis.

References

- Guida, G. & Mauri, G. (1984). A Formal Basis for Performance Evaluation of Natural Language Understanding Systems. *Computational Linguistics*, 10, 15-30.
- Guida, G. & Mauri, G. (1986). Evaluation of Natural Language Processing Systems: Issues and Approaches. *Proceedings of the IEEE*, 74, 1026-1035.
- Woods, W. A., (1977). A Personal View of Natural Language Understanding. *SIGART Newsletter*, 17-20.

A Sample Exemplar

- (1) The next day after we sold our car, the buyer returned and wanted his money back. (Allen, 1987, p. 346)
- (2) The day after we sold our house, the escrow company went bankrupt.
- (3) The day after we sold our house, they put in a traffic light at the corner.

Topic

Anaphoric reference - roles.

Discussion

In (1) the 'buyer' refers back to a figure in one of the roles in the 'selling a car' event. The system must search not only the direct possible antecedents (the 'selling') but must also consider aspects of the selling to resolve the reference. In (1), there is nothing specific to 'car' about resolving the reference. But in (2), finding the reference of 'the escrow company' involves looking past the general "buying" script and searching through aspects of selling specific to selling houses. There is a general problem here with controlling the amount of search while still looking deep enough. In (3), the system has to go from the house to the location to the street to the corner to understand the reference.

Reference

Allen, J. F. (1987). *Natural Language Understanding*. Menlo Park, CA: Benjamin/Cummings.